

multi-head Attention in Transformers

what is the drawback with
the self-attention?

there is ambiguity
in the sentence.

[The man saw the astronomer with a telescope]

	The	man	saw	the	astronomer	with	a	telescope
The								
man								
saw								
the								
astronomer								
with								
a								
telescope								

There are 2m coming to this

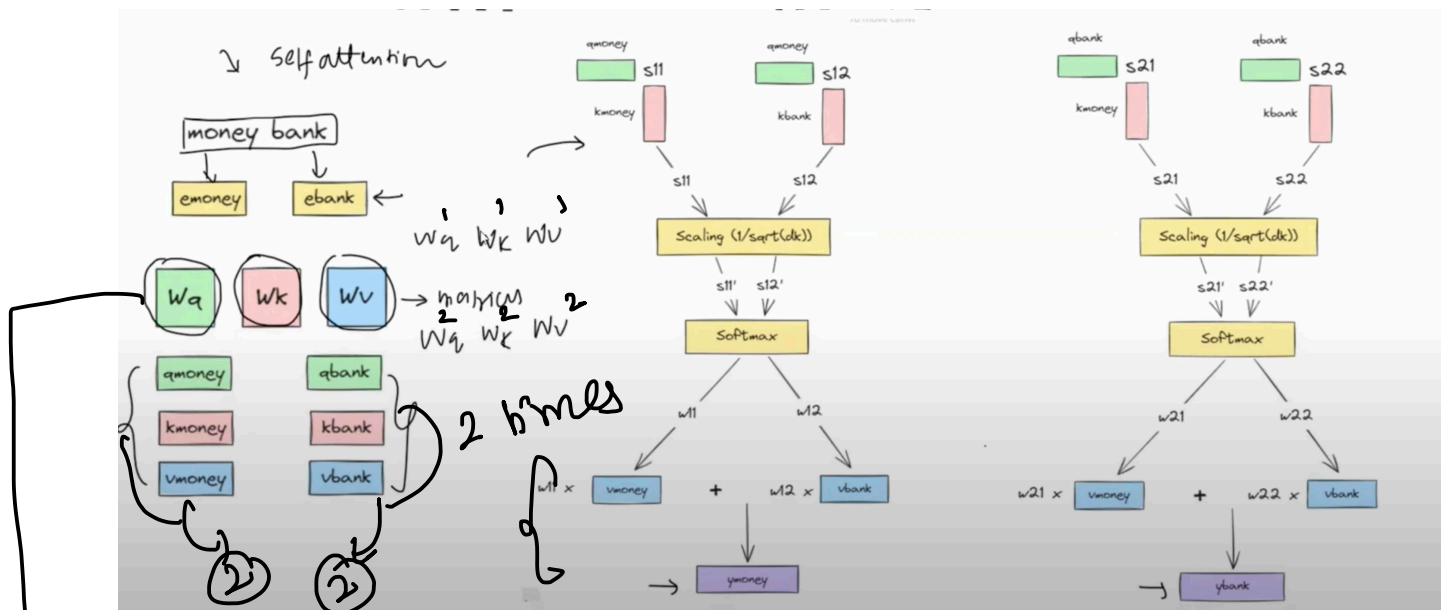


telescope

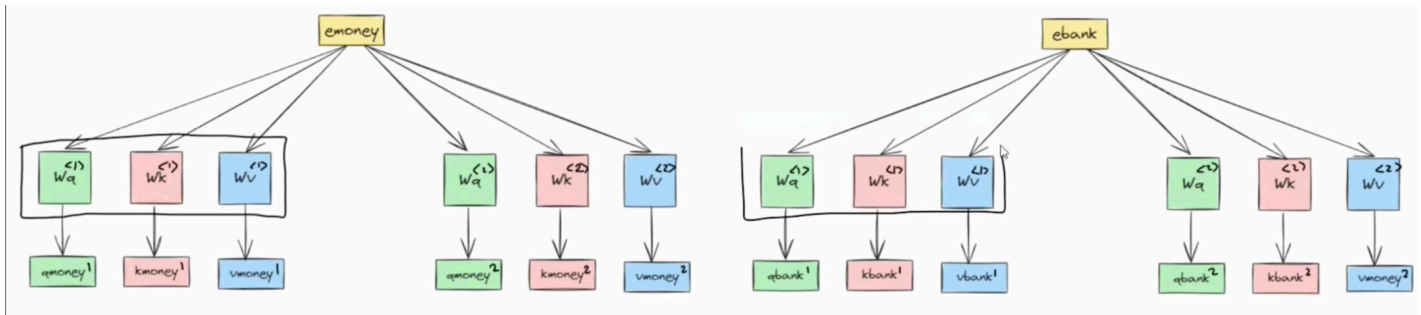
∴ But self-attention can capture only one meaning.

To solve this issue, multi-head attention was introduced.

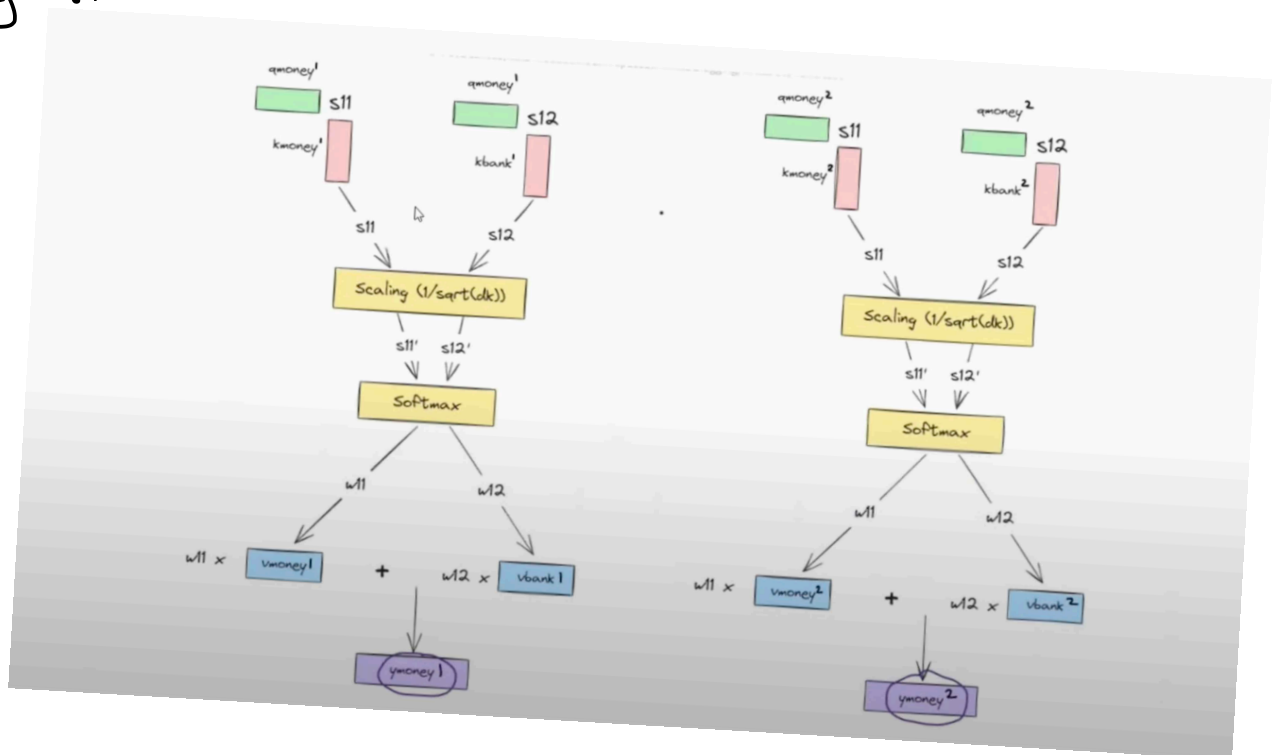
∴ Basically multi-head attention is more than one self-attention. One head is one self-attention.

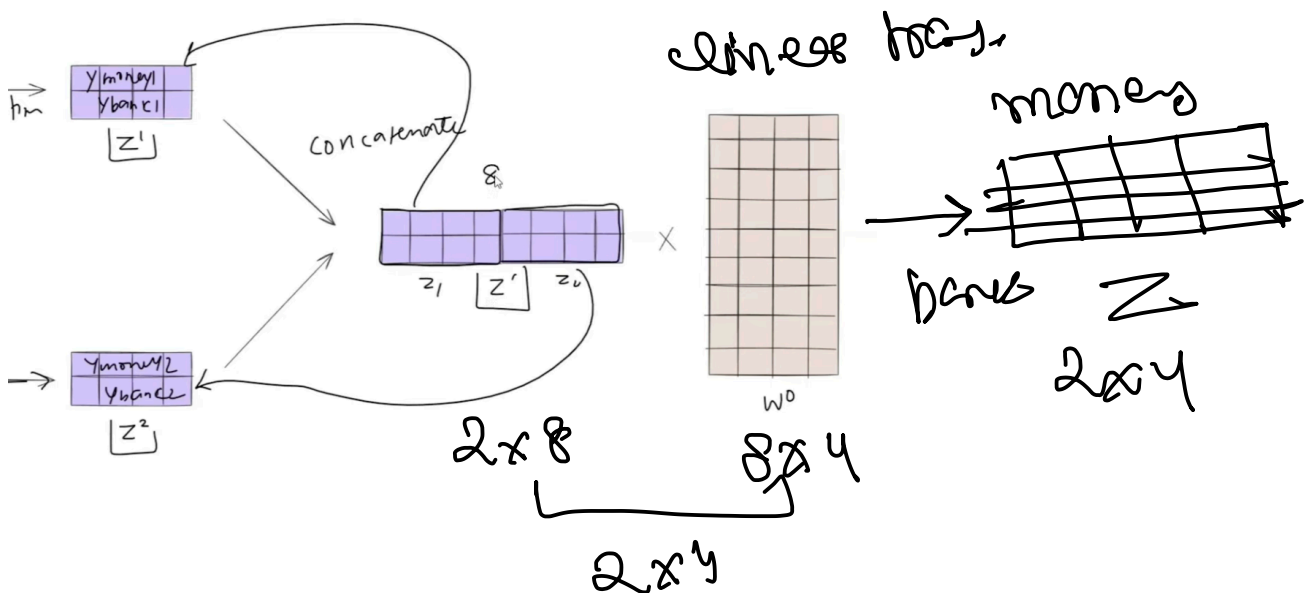
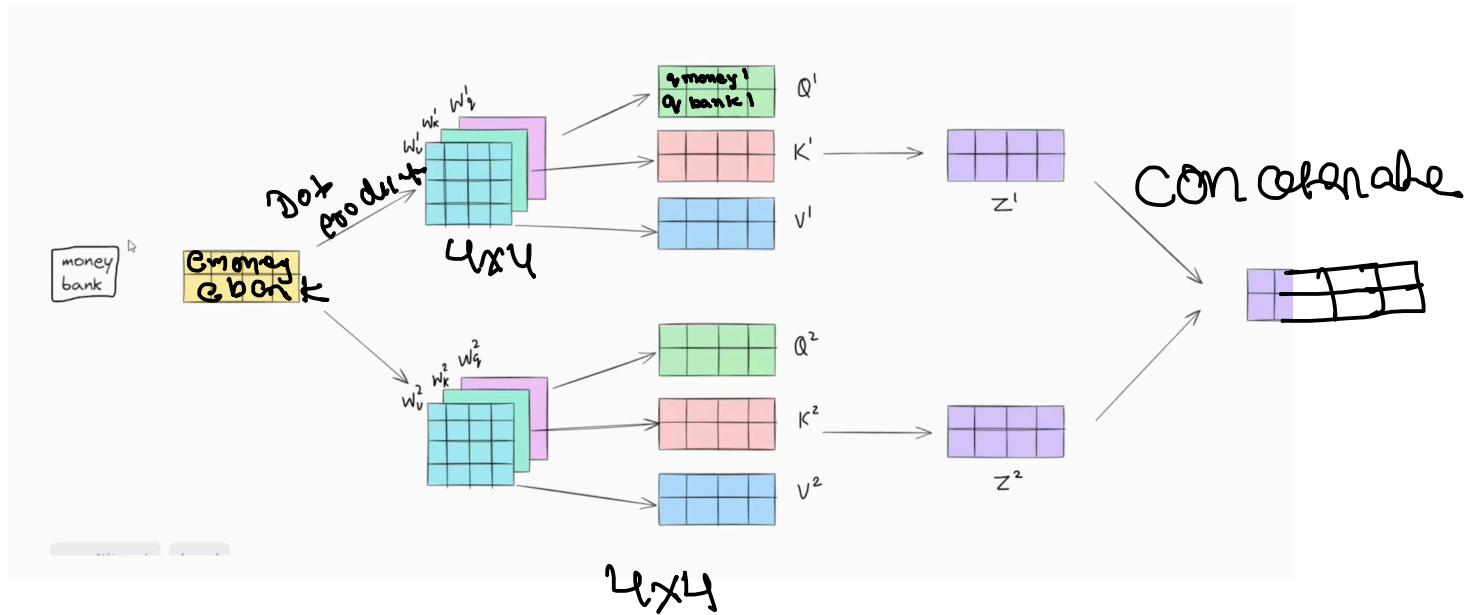
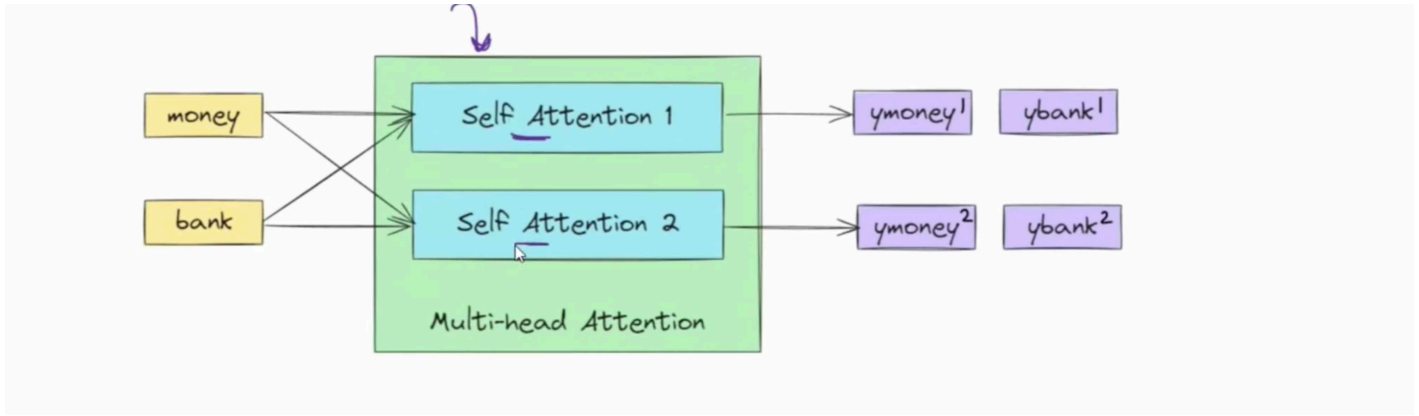


When working on self-attention, we have only one set of w_q , w_k , w_v . But in case of multi-head attention, we have this set more than once.

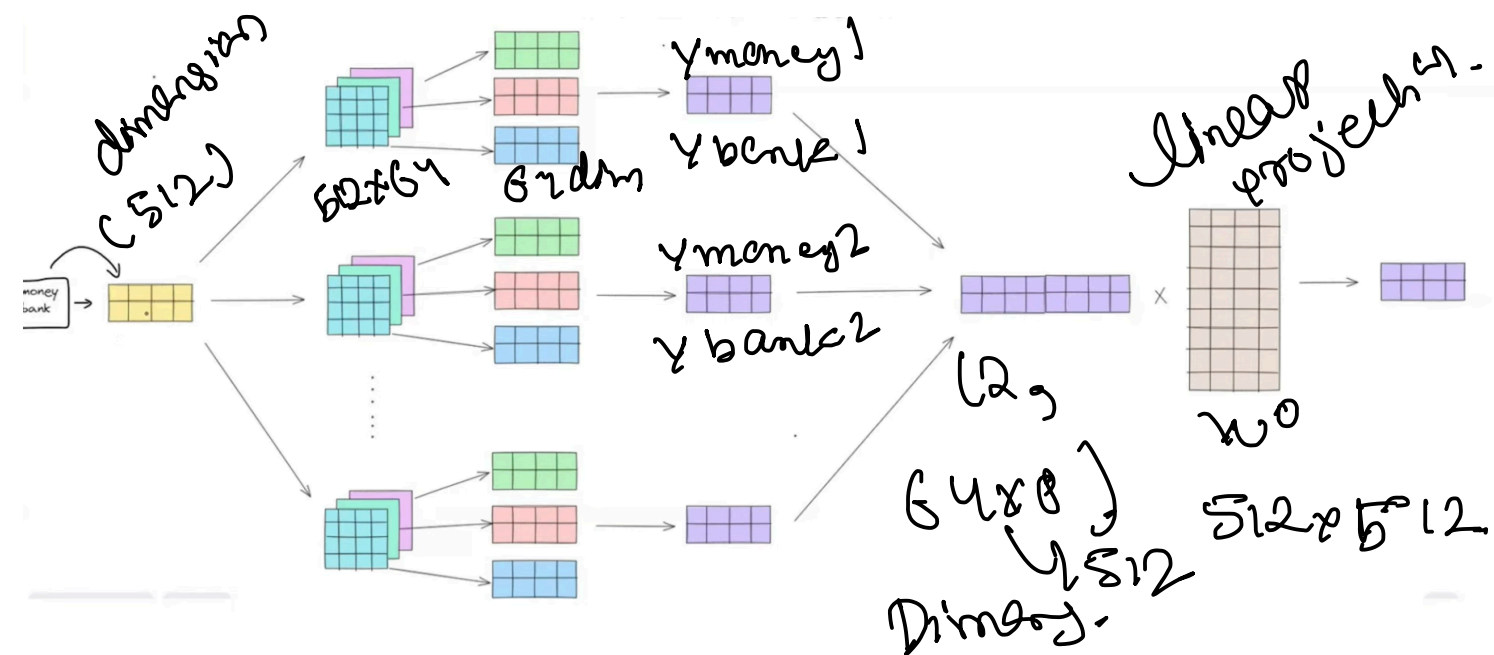


As we see below, 2 **vmoney** are generated.





Attention is all you need
 & attention



we reduce
 the Dimension
 Because we wanted to
 reduce the computation.